

Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/122998/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Jagadeesh, Karthik A., Paggi, Joseph M., Ye, James S., Stenson, Peter D., Cooper, David N. ORCID: <https://orcid.org/0000-0002-8943-8484>, Bernstein, Jonathan A. and Bejerano, Gill 2019. S-CAP extends pathogenicity prediction to genetic variants that affect RNA splicing. *Nature Genetics* 51 (4) , pp. 755-763. 10.1038/s41588-019-0348-4 file

Publishers page: <http://dx.doi.org/10.1038/s41588-019-0348-4>
<<http://dx.doi.org/10.1038/s41588-019-0348-4>>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies.

See

<http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



Title

S-CAP extends clinical-grade pathogenicity prediction to genetic variants that affect mRNA splicing

Authors

Karthik A. Jagadeesh^{*1}, Joseph M. Paggi^{*1}, James S. Ye¹, Peter D. Stenson², David N. Cooper², Gill Bejerano^{1,3,4,5}

Affiliations

¹ Department of Computer Science, Stanford University, Stanford, California 94305, USA

² Institute of Medical Genetics, Cardiff University, Heath Park, Cardiff, UK

³ Department of Pediatrics, Stanford University, Stanford, California 94305, USA

⁴ Department of Developmental Biology, Stanford University, Stanford, California 94305, USA

⁵ Corresponding Author: bejerano@stanford.edu

* denotes equal contribution

Abstract

There are over 15,000 known variants that cause human inherited disease by disrupting splicing [GILL: According to the latest version HGMD, we have logged ~19,800 splicing-relevant lesions out of a total of ~220,000 mutations (9%)]. While there are a few *in silico* methods, such as CADD, EIGEN and LINSIGHT, that predict the pathogenicity of noncoding variants, none are focused explicitly [GILL: Do you mean exclusively?] on splicing and none are able to effectively distinguish pathogenic splicing variants from benign variants. We introduce S-CAP, a novel splicing pathogenicity predictor that reduces the number of splicing variants of uncertain significance in patient exomes by 41%, a nearly 3-fold improvement over existing noncoding pathogenicity measures, while correctly classifying known pathogenic splicing variants with a clinical-grade 95% sensitivity.

Introduction

Genomic sequencing, and in particular exome sequencing, is revolutionizing the diagnosis of Mendelian disease, with over 5,000 genetic diseases already successfully mapped to over 3,000 genes [GILL: I think you will find that this is an underestimate, both in terms of the number of genetic diseases and the number of disease genes]. The sifting of patient exomes in search of a causal variant is a time-consuming process that often focuses on the coding sequence (CDS) of genes¹⁻⁴. Powerful pathogenicity meta-predictors such as M-CAP⁵ integrate multiple primary predictors such as SIFT⁶ and Polyphen-2⁷ with cross-species sequence conservation features to offer accurate clinical grade predictions capable of missing only a tiny fraction [GILL: a small proportion?] of known coding pathogenic mutations. This efficiency is important [GILL: you could say that both sensitivity and specificity are important?], because reducing the size of the candidate list of variants of uncertain significance (VUS) is futile if the pathogenic variant itself is incorrectly classified as benign⁸.

M-CAP is commonly used by clinicians to prioritize nonsynonymous variants^{9,10}, so that they can effectively consider first the variants that are most likely to yield a diagnosis. However, with diagnostic rates at 30-50%, one must look beyond the CDS itself. Splicing is a complex and crucial step of gene expression, wherein vast sections of RNA, called introns, are removed from a pre-messenger RNA and the remaining RNA, called exons, are joined together to form the mature messenger RNA (mRNA). Changes in splicing induced by genetic variants can have severe impacts on the protein coding potential of an mRNA, such as the exclusion of an entire exon¹¹ or a frameshift^{12,13} induced by creation of a new splice site¹⁴, among other effects¹⁵⁻¹⁷. Exome sequencing captures sequence information up to 50 base pairs past exon boundaries into each adjacent intron¹⁸. This region covers a broad class of splice affecting [GILL: splicing-relevant? splice-disrupting?] variants: those that disrupt existing splice sites or exonic and proximal intronic splicing regulators, such as the branchpoint, and some that create new splice sites.

Indeed, there are nearly 20,000 known Mendelian disease causing variants that impact the gene product through RNA splicing^{19,20}. However, a typical singleton patient's exome contains over 500 variants of uncertain significance (VUS)²¹ and (as we show) there currently exist no tools adequate for the purpose of interpreting these variants. Pathogenicity prediction tools such as CADD²², EIGEN²³ and LINSIGHT²⁴ have attempted to tackle a broad spectrum of non-coding mutations, but in doing so they neither focus on predicting splicing variant pathogenicity, nor do they provide any clinical grade assurances of minimizing the false prediction of known pathogenic splicing variants as benign. Generic methods also ignore the rich literature characterizing mRNA splicing and predicting associated molecular phenotypes, such as the percentage-spliced-in of exons²⁵⁻²⁸. These findings and methods provide invaluable insight into the potential of a variant to disrupt splicing. However, the splicing literature does not tell the whole story either, as predicting molecular phenotypes is a fundamentally different task from predicting if a variant will cause a disease. For a variant to cause a disease, the variant must disrupt normal splicing in one or more relevant tissues *and* the induced change in mRNA phenotype must be pathogenic. Similar obstacles are met in the few cases where clinicians attempt to identify patient genetic variants that disrupt splicing by performing a costly and time consuming RNA-seq experiment in an accessible, but not necessarily disease relevant, cell population^{29,30}.

We introduce S-CAP (**Splicing Clinically Applicable Pathogenicity**), a machine learning tool that integrates knowledge of splicing with measures of variant, exon and gene importance into a splicing-specific pathogenicity score. We evaluate S-CAP at the high sensitivity required in clinical settings and show that it far outperforms existing non-coding pathogenicity scores, as well as tools focused solely on identifying synonymous variants that disrupt splicing. S-CAP will allow clinicians to consider a broad class of splice-impacting variants, resulting in the diagnosis of more patients suffering from Mendelian disease.

Results

We developed a machine-learning framework to model and evaluate splicing variant pathogenicity. We analyzed the positional distribution and potential functional effects of variants near splice sites and defined 6 regions that displayed distinct mutation rates and functional effects. We trained 6 models (1 per region) to predict splicing variant pathogenicity. This involved building (i) a labeled dataset of known pathogenic and benign splicing variants, (ii) a set of features to help discriminate between the pathogenic and benign variants in each region and (iii) a learning algorithm to identify patterns in the features and to distinguish between variants in the two classes. Finally, we evaluated our models on a [[GILL: presumably unrelated] set of known pathogenic and benign variants, as well as over a dozen individuals with diseases caused by splice affecting [GILL: splicing-relevant? splice-disrupting?] variants.

The landscape of splice region variation

To build a dataset of labeled splicing variants, we considered variant pathogenicity, semantic effect and population frequency. We start by taking the union of 109,279 pathogenic single nucleotide variants (SNVs) from the Human Gene Mutation Database (HGMD) and 25,793 pathogenic SNVs from ClinVar to obtain a total of 114,382 unique pathogenic variants [GILL: you gained ~5,000 extra variants (4%) by including data from ClinVar. Personally, I am a little doubtful that this exercise was worth it. I may be somewhat biased but I do wonder if the additional ClinVar data are going to be as reliable as the mutation data from HGMD which have (i) all been peer-reviewed prior to publication and (ii) been independently assessed by our expert curators prior to inclusion in HGMD. I would not be surprised if your AUC values improved slightly if you confined your analysis to HGMD data. Maybe you should try this? No harm done if I am wrong....and if I am right, your results will be even more impressive]. We curated 15,833,389 benign SNVs observed in controls from the gnomAD database who do not suffer from any obvious Mendelian disease. To identify a subset of variants with likely impact on splicing, these sets were filtered to the 'splicing region', i.e. all synonymous or intronic variants within 50 base-pairs^{18,31} of an exon boundary (we justify the choice of name [GILL: term? splicing region?] below). Removing nonsynonymous and loss of function (stop gain or stop loss) variants ensured that our model was trained and evaluated virtually exclusively on splicing variants. The gnomAD variants were further filtered to remove any variant identified as being pathogenic in HGMD or ClinVar. This resulted in 14,938 splicing pathogenic variants and 7,027,609 splicing benign variants. Then a frequency filter, based on the ACMG guidelines that suggest clinicians consider common (> 1% frequency) variants as definitively benign, was applied to both sets yielding 14,838 rare splicing pathogenic variants and 6,760,450 rare splicing benign variants. In support of this ACMG guideline, we note that only 100 of 14,938 (0.67%) known pathogenic variants in the splicing region are common in the general population (see Methods).

To show that nearly all pathogenic 'splicing region' variants do in fact disrupt splicing, we developed a simple model to assign a putative effect of variants on splicing. Among other factors, splice affecting [GILL: splicing-relevant? splice-disrupting?] mutations can (1) create a cryptic splice site (2) disrupt an existing splice site or (3) disrupt an existing branchpoint. We used the existing tools MaxEntScan and LaBranchoR, which predict the strength of splice site sequences and branchpoint sequences, respectively. We denoted a variant as creating a cryptic splice site if the variant creates a splice site with a MaxEntScan score at least as high as the score for the reference splice site (with the variant included), disrupting a splice site if it reduces the MaxEntScan³² score of the reference splice site by more than 1, and as disrupting a branchpoint if it has a LaBranchoR³³ *in silico* mutagenesis score of less than -0.1. In cases where a variant had multiple putative effects, e.g. creating a cryptic splice site and disrupting an existing splice site, we resolved to the most extreme effect, which we took to be the order in which they are introduced above [GILL: I am unclear as to your meaning and even more unclear in relation to how you establish the order in which 'they are introduced'. Do you mean the order in which they are likely to impact the RNA splicing phenotype?]. We found that 97% of putatively pathogenic splicing variants are predicted to have an effect on splicing, as opposed to only 18% of likely benign variants (**Fig. 1a-b**). We also found that pathogenic variants are enriched at positions where the mechanistically essential U2 snRNA, U2AF, and U1 snRNA bind. Further, the positional distribution of pathogenic variants traces the bias in nucleotide content inside of these binding sites [GILL: meaning unclear!](**Fig. 1b**). Conversely, variants occurring in the general population are biased away from these high information content positions³⁴ (**Fig. 1a**).

Region-specific models to increase performance and alleviate ascertainment bias

In order to effectively capture these position-specific patterns, we separated the variants into 6 regions relative to the splice sites with generally homogenous function and built a separate model for each region (**Supplementary Fig. 1**). Specifically, we grouped variants occurring in the obligate 5' GT (5' core) and 3' AG (3' core) dinucleotides, intronic variants upstream of a 3'ss (3' intronic), variants lying in the canonical U1 snRNA binding site excluding the core 5'ss (5' extended), intronic variants downstream of a 5'ss (5' intronic), and synonymous variants within the protein coding gene (exonic) (**Fig. 1c**).

Core splice site variants are well known as having a large functional effect and are easy to identify, a fact that has probably led to an overrepresentation of core splice site variants in pathogenic variant databases. Around 73% of known splicing region pathogenic variants occur within the core splice sites (**Fig. 1c**). Although this is consistent with the mechanistic importance of these positions, in unbiased studies of splicing quantitative trait loci it is generally found that fewer than 1% of splicing affecting [GILL: splicing-relevant? splice-disrupting?] variants are found to be located in core splice sites^{35,36} and in a recent study employing RNA-seq data to identify splice-affecting [GILL: splicing-relevant? splice-disrupting?] variants resulting in Mendelian disease, only 2 of 6 (33%) causal variants in the splicing region were in core splice sites²⁹. If left unaddressed, this bias would allow a classifier to have strong test set performance (by calling most core splice site variants pathogenic and others benign), but would often miss non-core splice site pathogenic variants. Separating variants by position allows us to guarantee that pathogenic variants are rarely misclassified as benign in every region, thereby ensuring an overall low false negative rate irrespective of ascertainment biases present in annotated data.

S-CAP features

We curated existing metrics and developed several novel features to help distinguish between pathogenic and benign variants within the splicing region (**Supplementary Table 1**). The set consists of chromosome, gene, exon and variant level features. At the **chromosome level**, we add 3 features to distinguish between variants found on chromosome X, chromosome Y and the autosomes. Variants on the X chromosome present an important subcase [GILL: subset?] since in males, a hemizygous X chromosome variant inducing loss of function results in no viable gene product. Consistent with this intuition, pathogenic variants are highly enriched on the X chromosome as compared to the autosomes (7.11 fold enrichment, $p < 10^{-140}$ by two-sided Fisher's Exact Test). At the **gene level**, pLI³⁴, RVIS³⁷, and a haploinsufficiency score³⁸ help to measure the potential pathogenicity of each gene [GILL: meaning unclear. Are you referring to the Human Gene Damage Index? Do you mean the likelihood that a given gene is a 'disease gene'?]. At the **exon level**, exon length, exon length modulo 3, reference splice site strengths, an existing regional constraint score³⁹ and the exon sequence similarity between hg19 and 99 species from the 100way alignment serve to assist in distinguishing critical exons from those that may be safely excluded [GILL: Have I understood this correctly?]. Additionally, we developed a novel splice site constraint score to measure the fragility and tolerance of each exon to splice site mutations (see Online Methods). At the **base-pair level**, CADD²² measures pathogenicity based on functional data annotations, LINSIGHT measures variants' fitness effect through functional data and molecular evolution, whereas SPIDEX²⁷ was incorporated in order to measure the impact of a given variant on exon inclusion. PhyloP⁴⁰ and PhastCons⁴¹ scores from the multiz46way and multiz100way alignments measure the evolutionary importance of the affected base across primate species, placental mammals and all vertebrates. We also include a feature to capture the change in 3-mer content induced by a variant⁴². Additionally, we include region-specific features, such as a branchpoint disruption term for the 3' intronic region³³ and a 5' cryptic splice site creation term for the 5' intronic region (see Online Methods for a complete description of features).

The machine learning algorithm

Similar to the M-CAP classifier⁵, S-CAP is built using a gradient boosting tree classifier, a highly effective machine learning model⁴³. This model iteratively builds decision trees, where each tree is picked to correct the most cases that were misclassified in the previous step. The final classifier is a linear combination of each of the previously derived decision trees (see Online Methods for details).

S-CAP consistently outperforms existing pathogenicity scores

Each of the 6 previously described regions contains a set of pathogenic and benign variants (**Fig. 1c**), which were used to train 6 separate models. We performed 5-fold cross-validation and selected the median performing model as the final model for each region. S-CAP was evaluated against the most popular existing methods that score splice affecting [GILL: splicing-relevant? splice-disrupting?] variants: CADD, SPIDEX, LINSIGHT and EIGEN. Existing methods all performed at random or better (AUC of 0.46 - 0.81) across the 6 defined splicing regions (**Fig. 2a-f**). The 6 S-CAP models outperformed all existing methods in all regions resulting in up to a 26.6% improvement in the AUC over the next best performing model. S-CAP performance ranged from achieving an AUC of 0.804 in the 3' core region (**Fig. 2b**) to achieving an AUC of 0.953 in the 3' intronic region (**Fig. 2a**). No existing method consistently outperforms other existing metrics across all splicing regions (**Fig. 2a-f**).

Additionally, S-CAP performance on exonic variants was independently compared against MutPred Splice (see Supplementary Figure 2)⁴², a tool focused on scoring only splice affecting [GILL: splicing-relevant? splice-disrupting?] synonymous variants. As MutPred Splice was trained using the same data used to build S-CAP, a random train and test split of data may have resulted in the inclusion of variants in the test set that were used to train the MutPred Splice classifier. To ensure zero information leakage between the training and test datasets, we carefully built a test set that excluded all MutPred Splice training data.

Clinically relevant threshold maintains high sensitivity

As previously shown in M-CAP⁵, it is important to tune thresholds for clinical settings so that fewer than 5% of known pathogenic variants are misclassified. None of LINSIGHT, EIGEN or SPIDEX provides a default threshold to consider a variant as pathogenic. As a result, it is difficult to use any of these methods for variant pathogenicity classification. CADD provides a threshold but, at the author-recommended default threshold, over 31% of the known pathogenic variants across the splice region are discarded and incorrectly classified as benign (**Table 1**). Similarly, MutPred Splice provides a default threshold, but over 48% of known pathogenic exonic variants are misclassified at this threshold. We generated a high sensitivity threshold for all metrics in each region by finding the lowest threshold that results in the correct classification of at least 95% of a test set of pathogenic variants from that region (**Table 2**).

S-CAP excels at clinically relevant thresholds

Moving into the high sensitivity domain, S-CAP's performance (**Supplementary Fig. 3a-f**) ranged from achieving an hsr-AUC of 0.186 in the exonic region (**Supplementary Fig. 3c**) to an hsr-AUC of 0.549 hsr-AUC in the 3' intronic region (**Supplementary Fig. 3a**). When evaluating MutPred Splice on an independent test set of exonic variants, S-CAP achieves an hsr-AUC of 0.204 and outperforms MutPred

Splice which achieves an hsr-AUC of 0.081. Overall, S-CAP improves on existing metrics by up to 472% when focused on the high sensitivity domain.

Different patterns observed for recessive and dominant variants

Recessive and dominant diseases are associated with different selective pressure on alleles, which we hypothesized would result in different feature importances and thresholds for determining variant pathogenicity. To address this complexity, we developed separate classifiers for dominant (heterozygous in patient) and recessive (homozygous in patient) alleles. In patients, we are given whether each variant appears in the heterozygous or homozygous state. However, our training data did not provide dominant and recessive labels, so we developed a framework for inferring this information from a control population for model training (see Methods).

Tagging 3' and 5' core variants as dominant or recessive resulted in improved performance **Supplementary Fig. 4a-h**) compared to models without these tags (**Fig. 2b,d**). For 3' core variants, we built a model that improved upon an AUC of 0.804 and hsr-AUC of 0.257 from the original 3' Core model to an AUC of 0.805 and 0.890 (**Supplementary Fig. 4c,d**) and hsr-AUC of 0.296 and 0.454 (**Supplementary Fig. 4g,h**) when tested on just dominant or recessive variants, respectively. Similarly, for 5' core variants, we built a model that improved upon an AUC of 0.805 and hsr-AUC of 0.291 from the original 5' core model to achieve an AUC of 0.779 and 0.880 (**Supplementary Fig. 4a,b**) and hsr-AUC of 0.222 and 0.518 (**Supplementary Fig. 4e,f**) when testing on just dominant or recessive variants, respectively. The S-CAP model to be used on patients takes advantage of these split dominant and recessive models.

S-CAP eliminates the most VUS in patient exomes

Resources like S-CAP are developed on large sets of benign and pathogenic variants but ultimately are used to help with the interpretation of VUS in individual patients (**Fig. 1C**). To demonstrate the practical utility of S-CAP, we evaluated S-CAP and each of the comparison methods on 14 patients with Mendelian diseases caused by splice-altering mutations. After applying the standard allele frequency filter of $\leq 1\%$, a typical individual has on average a total of 533 rare variants within the splicing region (**Fig. 1C**). Typically, $\sim 32\%$ of variants were observed in each of the 5' intronic, exonic and 3' intronic regions. 3% were in the 5' extended regions and 0-3 variants were in each of the 3' core and 5' core regions (**Fig. 1c**). Existing methods, CADD, LINSIGHT, EIGEN and SPIDEX, using the 95% sensitivity thresholds (**Table 2**), perform comparably when applied to all VUS in the splicing region of an individual patient and on average reduce the number of VUS by up to 15%. By contrast, S-CAP is much more powerful, reducing the number of VUS in the splicing region for an individual patient by 31-46% (**Fig. 3c**), while confidently retaining the pathogenic variant for further detailed analysis (**Table 3**). For a larger sample size, we evaluated each method on all ($n=2054$) individuals in the 1000 Genomes Project²¹, which can conceptually be thought of as Mendelian disease patients with their pathogenic variants removed (see **Fig. 1C**). The observed fraction of VUS reduced on average per 1000GP individual is consistent with the performance observed when applied to patients (**Fig. 3d**). Specifically, S-CAP is nearly three times as powerful as existing methods and on average results in a 41% reduction of the VUS within the splicing region of a given individual.

Discussion

Variants affecting splicing comprise the second largest category of known pathogenic mutations²⁰. A broad class of potential splice affecting [GILL: splicing-relevant? splice-disrupting?] variants are already being captured by exome sequencing, yet clinicians often ignore these variants because they do not have the proper tools to interpret them. Here we address this problem by developing S-CAP, the first clinically applicable pathogenicity predictor dedicated exclusively to splicing variants.

In order for an *in silico* pathogenicity predictor to be useful in the clinic, it needs to be easy to use, carefully evaluated at high sensitivity and confidently removes a substantial fraction of benign variants. Existing noncoding variant tools are not easy to use for splicing variants because their performance is strongly dependent upon position relative to splice sites, classification thresholds are either not provided or poorly calibrated, and no single method consistently outperforms the others. This means that employing existing methods, a clinician would have to consult multiple pathogenicity scores for variants in different regions. Furthermore, none of the existing methods have been carefully evaluated at clinical-grade sensitivity and most methods give no guidance about what cutoff should be used. Of the methods that do suggest a cutoff, all result in the misclassification of an unacceptably high number of pathogenic variants (**Table 1**). After carefully evaluating the existing methods, we found that none of them confidently removed a considerable fraction of benign variants (**Figure 3**). For example, after we retuned it (**Table 2**), SPIDEX performed well on variants in the intronic bins, but was close to random at predicting the pathogenicity of core splice site variants. After retuning, CADD (**Table 1,2**) performed well on core and exonic variants but poorly on intronic variants (**Figure 2**). S-CAP, addresses these important issues, as it consistently outperforms existing methods across all the regions, its performance has been carefully evaluated at clinical-grade sensitivity, and it removes close to three times as many benign variants as any other existing method (**Figure 3**).

Central to the design of S-CAP is the use of region-specific models to alleviate the effects of ascertainment biases in curated pathogenic variant databases. Curated pathogenic variant databases contain invaluable information about the properties of pathogenic variants, but they also over-represent variants in known disease genes and [GILL: in association with?] easily identifiable features. Of particular concern in splicing pathogenicity prediction is the inflated number of variants in core splice sites, which exists because they are easily recognized and have well established molecular consequences. If left unaddressed, this bias in the labeled pathogenic data would lead to unrealistic model performance, as a model could achieve relatively high test set performance simply by predicting that all core splice site variants are pathogenic and all others are benign. However, in the clinic, such a model would incorrectly classify pathogenic, non-core splice site variants as benign at an unacceptably high rate. Introducing separate models for each region alleviates this concern since each model's performance is evaluated using data from the same region, thereby assuring high sensitivity irrespective of the underlying positional distribution of pathogenic variants. Additionally, we allowed for the over-representation of pathogenic variants in known disease-associated genes by ensuring that variants from the same gene were never split between the training and evaluation sets (see Methods). This guaranteed that no gene level information was shared by features across folds.

The use of patient RNA-sequencing (RNA-seq) data to identify pathogenic variants that disrupt splicing is a growing and promising field^{29,30,44} to which we believe S-CAP is complementary. In fact, our model already includes scores from SPIDEX²⁷, a deep learning model that was trained on tissue-specific RNA-seq data to predict the change in percent spliced-in ($\Delta\Psi$) of an exon given a variant. It would only be a small step to supplementing these predicted $\Delta\Psi$ s with experimentally measured $\Delta\Psi$ s from RNA-seq experiments. Observing $\Delta\Psi$ through RNA-seq bypasses the difficult problem of explicitly or implicitly predicting the effect of a given variant on splicing in a particular cell type, but this is only part of the problem. Whether predicted or measured, gene expression and Ψ values vary between cell contexts and time points and in many cases the most relevant cell population to sequence would not be clear^{30,45}. Perhaps even more importantly, molecular phenotypes are multiple steps away from real phenotypic

change, making it difficult to predict organismal pathogenicity from molecular phenotype⁴⁶. These factors limit the direct applicability of observed or predicted $\Delta\Psi$ s to cases where there is a good understanding of the relationship between a disease, the cell population it affects, and the set of potentially causative genes. Many of the features used by S-CAP, such as evolutionary conservation, implicitly integrate an allele's importance over all cell populations and time points, complementing the strengths of RNA-seq based methods. Another direction for the integration of S-CAP and experimental methods is the use of cheap and fast site-directed sequencing in a diverse array of cell types to validate putatively pathogenic sites identified by S-CAP⁴⁷.

Eventually, analyzing the splicing region will become commonplace in clinical settings. This is currently a difficult task given the complexity of splicing and the difficulty in predicting whether a change in splicing will result in disease. There are over five hundred rare variants of uncertain significance per individual in the splicing region with no clear semantic effect outside of the core splice sites and most are not observed in control populations. S-CAP represents a big step towards effectively interpreting splicing variation, but we will have to continue to improve on these methods, learning more from RNA-seq experiments, to render this problem more tractable.

Online Methods

Variant Processing

Dataset of pathogenic and benign variants

Pathogenic variants were obtained from two manually curated databases: the Human Gene Mutation Database (HGMD) Professional version 2017.1 and ClinVar release 20170406. Only HGMD variants tagged as Disease Mutation (DM) and ClinVar variants with Pathogenic Clinical Significance were included in the final set of pathogenic variants [GILL: How many variants from HGMD, how many from ClinVar?]. Benign variants were obtained by identifying variants observed in individuals from gnomAD³⁴ r2.0.2.

Variant Annotation

ANNOVAR⁴⁸ v527 was used to annotate variants with predicted effect on protein-coding genes using gene isoforms from Ensembl⁴⁹ gene set version 75 for the hg19/GRCh37 assembly of the human genome. All coding isoforms were used where the transcript start and end sites were marked as complete and the coding span was a multiple of three.

Variant Filtering

All variants were filtered so as to only include rare variants that do not directly affect the protein coding sequence and the splicing region [GILL: Meaning unclear. Variants outwith the obligate GT and AG dinucleotides or extended splice site consensus sequences? Same definition as given in Results i.e. all synonymous or intronic variants within 50 base-pairs^{18,31} of an exon boundary?]. Rare variants are defined to be variants with an allele frequency of $\leq 1\%$ in all control populations and subpopulations in KGP phase 3, ExAC v0.3.1 and gnomAD r2.0.2. Variants that do not affect the protein coding sequence are those determined to be synonymous or which in the core splicing, extended splicing or intronic regions.

Positional Subsetting of Variants

Variants at different positions relative to splice sites have different properties, such as distributions of sequence conservation and ratio of pathogenic variants to benign variants. This led us to allocate

variants into subsets based upon their position relative to splice sites and train separate models for each. Subsets were selected so as to group together functionally related positions. In total, we constructed 6 subsets (Fig. 1). We created subsets of 3' and 5' core (1 or 2 bases from the splice site), extended (1-2 bases into the exon on the 5' side and 3 to 6 bases from the 5' splice site), 3' intronic (2 to 50 bases from the splice site), 5' intronic (7-50 bases from the 5' splice site) and exonic (inside the exon but outside the extended region) regions.

We associated each variant with a nearby [GILL: neighboring?] exon. For variants close to multiple exons, we attempted to assign the variant to the exon considered as having the highest chance of a pathogenic effect, which we inferred from the density of pathogenic variants in each bin. Specifically, we favored associations in the following order: core, 5' extended, intronic, and exonic.

Features

Our models utilized a diverse library of previously described and novel features (fig. S1). These can be divided into gene level, exon level, and variant level features as well as a few features that are specific to individual regions. Below, we introduce our novel features and describe how we curated previously described features.

Gene Level

We obtained **RVIS**³⁷ scores from, the appropriately named, genic-intolerance.org (see URLs). We used the data in column "RVIS[pop_maf_0.05%(any)]" as a feature in our models. We obtained **pLI**³⁴ scores from supplementary table 13 of the original publication on 7 April 2017. We used as a feature the column denoted as pLI. We obtained a **haploinsufficiency score**³⁸ which is a probability of a gene being haploinsufficient directly from the publication page on the PLOS website in May 2017. We downloaded MPC, a recently proposed regional constraint score³⁹ (see URLs).

Exon Level

For each exon, we created **splice site features** which measure the number of rare and common variants observed in gnomAD in the 5' and 3' core (0-2) and extended (2-6) regions. This was motivated by the desire to share information between functionally similar positions. We take this feature to represent the constraint on the exon's splicing region in the human population. To avoid data leakage, when constructing this feature for a particular variant, we masked the effect of the variant itself on this score. Specifically, if a given 5'ss had one associated core variant, it was assigned a count of 0. If it had 2, both were assigned a count of 1. Additionally, we measured **exon identity** across vertebrates and found that the exon identities in many organisms were highly correlated. Principal components analysis (PCA) of the identity scores for all exons showed that 5 components explain the vast majority of variation in the data. To prevent overfitting, we included the original exon identities projected onto these first 5 principal components as features. We also included exon length and exon length mod 3 as features associated with each variant.

Variant Level

In order to consider the local sequence context of variants, we included **spectrum kernel features** representing the change in trinucleotide content induced by the variant⁴². For all 64 possible trinucleotides, we created a vector counting the number of occurrences in the alternative sequence and subtracted an equivalent vector for the reference sequence. The **MaxEntScan reference score**, **alternative score and difference** between the two were all used as features to quantify the strength of each exons' reference and alternative 5' and 3' splice site. **SPIDEX** scores²⁷ were downloaded directly from the Deep Genomics website (see URLs) on 7 April 2017. Any variant not assigned a SPIDEX score was assigned a value of 0. We included 8 scores measuring evolutionary conservation, specifically, **PhyloP 46way vertebrates**, **placental mammals**, **primates**, **PhyloP 100way vertebrates**, **PhastCons**

46way vertebrates, placental mammals, primates and the PhastCons 100way vertebrates. Any base not annotated with a conservation score was assigned a value of 0. We also included the signed **distance to the 5'ss and 3'ss splice site of the associated exon** as features to measure base-pair importance.

Region Specific features

The 3' intronic S-CAP model includes a **branchpoint feature** modeled using LaBranchoR, a bi-directional LSTM (long short-term memory) model trained on the genome sequence surrounding experimentally validated branchpoint sites³³. Specifically, we used the *in silico* mutagenesis scores available online (see URLs) as a feature.

We explicitly represented the strength of **cryptic sites** created by each variant using MaxEntScan³². In 3' intronic, 3' core, and exonic bins, we included a 3' cryptic splice site creation term and, in the exonic, 5' extended, 5' core, and 5' intronic bins, we included a 5' cryptic splice site creation term. For each variant, we scanned for the highest scoring splice site motif that overlaps the variant, excluding reference splice sites. We used as features the strength of the cryptic site, the change in the strength of the cryptic site induced by the variant, the distance from the reference splice site to the cryptic site and the difference in strength between the cryptic site and the reference splice site.

Model Training and Testing

We performed 5-fold cross-validation to train and identify a generalizable S-CAP model. 5-fold cross-validation refers to splitting the data into 5 roughly equally sized parts (folds). All variants found in a single gene were included in the same fold to ensure that there was no leakage of feature information across the training and test sets. We then merged 4 of 5 sets to form a training dataset, trained the model on this training dataset and evaluated on the remaining fold to obtain an expected accuracy. We performed this process 5 times (each combination of 4 folds was merged together to form the training dataset) testing on the fold that was not included in the training dataset.

To train the S-CAP model, we used a Gradient Boosting Tree model implemented in the python 2.7.13 sklearn version 0.18.1 library and used the default parameters to reduce the chance of overfitting the model. After training 5 models during the cross-validation phase, we picked the median performing model as the final classifier. The ROC curves were built based on performance on the test set for this specific median model.

Comparison Metrics

We sought to compare our performance to those of other methods used to infer the importance / pathogenicity of noncoding variants. We evaluated the performance of each of the methods below by using the output score directly as a pathogenicity score. For SPIDEX, we negated the score, as is consistent with large negative scores having a larger impact on function and constraint, respectively. We report performances for all subsets where the method reported scores for at least 50% of variants. Variants where a method did not report a score were excluded from evaluation.

CADD²² v1.3 scores were downloaded from the CADD website (see URLs). **SPIDEX**²⁷ scores were downloaded directly from the Deep Genomics website (see URLs) on 7 April 2017. Any variant not assigned a SPIDEX score was assigned a default value of 0. **LINSIGHT**²⁴ scores were downloaded directly from the LINSIGHT website (see URLs) on 7 April 2017. Any variants not assigned a LINSIGHT score was defaulted to be 0. **Eigen v1.0 coding** and **Eigen v1.0 noncoding**²³ Phred scores were downloaded from the Eigen website (see URLs). **MutPred Splice**⁴² score were downloaded from the MutPred Splice website (see URLs).

Recessive v. Dominant Classifiers

We developed separate classifiers for recessive and dominant acting variants in the 3' and 5' core splice site regions. We opted not to include dominant and recessive classifiers for the other 4 regions as we did not have a sufficient number of pathogenic variants to train and evaluate multiple models and, in our exploration, they made less of an impact. Intuitively, for core variants the molecular phenotype is obvious, a loss of splicing. This places the full burden on predicting if this change will result in disease, a task heavily dependent upon whether the variant acts via a recessive or dominant mechanism. Whereas in the other regions, the majority of possible variants have little impact on splicing, making predictions as to whether or not the variant will have an effect on splicing is the primary challenge, a task unrelated to inheritance mode.

When evaluating a patient, core variants observed as heterozygous are routed to the dominant classifier and variants observed to be homozygous are routed to the recessive classifier. Since there was no genotype information available for the labeled pathogenic and benign variants, we developed a framework to label variants as dominant or recessive based on their occurrence in a control population. We labeled pathogenic variants observed as heterozygous in the control population as likely recessive, because a single copy can be harbored with no major issues, and pathogenic variants never observed in the control population as likely dominant. Benign heterozygous variants that are never observed as homozygous in the healthy control population provide little information regarding their potential to cause a recessive acting disease. In this case, we tagged benign heterozygous variants never observed to be homozygous as dominant, and benign variants observed as homozygous as recessive. Additionally, as any variant on the X chromosome resembles a homozygous autosomal variant [GILL: I don't follow this reasoning. You mean in males or females?], all X chromosome variants were labelled as recessive [GILL: Seems reasonable].

We encoded whether each variant was considered dominant or recessive as a binary feature and trained a single gradient boosting tree model for each region. Then, we found two 95% true positive rate thresholds separately based on test sets of only dominant-tagged and only recessive-tagged variants.

Patient Datasets

Sequencing and diagnosis for all patients was performed by other laboratories. All patient data were downloaded by requesting access to the European Genome-Phenome Archive (EGA) and the database of Genotype and Phenotype (dbGaP) databases. Variant call files (VCFs) for 11 patients were submitted by the Deciphering Developmental Disorders (DDD) study in the European Genome-Phenome Archive (EGA) study EGAS00001000775. An additional 3 patient VCFs were submitted to the database of Genotype and Phenotype (dbGaP) study phs000655.v3.p1.

Data availability.

S-CAP scores for all rare variants in the predefined splicing region in the human genome, along with the source code, and final trained model for the S-CAP classifier, are available through the S-CAP website (see URLs), licensed under a Creative Commons Attribution-NonCommercial 4.0 International License. The S-CAP code repository is also available at Bitbucket (see URLs).

URLs

S-CAP website, <http://bejerano.stanford.edu/mcap>; S-CAP codebase [GILL: This URL does not connect with anything as yet], https://bitbucket.org/bejerano/splicing_classifier; RVIS, http://genic-intolerance.org/data/RVIS_Unpublished_ExACv2_March2017.txt; SPIDEX, <https://www.deepgenomics.com/spidex-noncommercial-download> [GILL: This URL doesn't work]; LINSIGHT, <http://compngen.cshl.edu/~yihuang/LINSIGHT/>; Haploinsufficiency, <https://doi.org/10.1371/journal.pgen.1001154> [GILL: This URL is for a manuscript, Huang et al. (2010) which is given as no. 38 in your list of references]; **LabranchOR**, <http://bejerano.stanford.edu/labbranchor/>; ExAC, ftp.broadinstitute.org/pub/ExAC_release/release1/regional_missense_constraint/

Acknowledgements

We thank Dr. Jon Bernstein and the members of the Bejerano Laboratory. P.D.S and D.N.C receive financial support from Qiagen through a license agreement with Cardiff University. This work was funded in part by Stanford Graduate Fellowship and the Computational and Evolutionary Genomics Fellowship to K.A.J. and the Stanford Pediatrics Department, DARPA, a Packard Foundation Fellowship and a Microsoft Faculty Fellowship to G.B.

Author Contributions

K.A.J., J.M.P. and G.B. designed the study. K.A.J., J.M.P. and J.S.Y. developed the features, trained the model and evaluated the results. P.D.S. and D.N.C. curated the HGMD data. K.A.J, J.M.P. and G.B. wrote the manuscript. All authors reviewed the manuscript.

Competing Financial Interests

The authors declare no competing financial interests.

References

1. Bamshad, M. J. *et al.* Exome sequencing as a tool for Mendelian disease gene discovery. *Nat. Rev. Genet.* **12**, 745–755 (2011).
2. Ng, S. B. *et al.* Targeted capture and massively parallel sequencing of 12 human exomes. *Nature* **461**, 272–276 (2009).
3. Ng, S. B. *et al.* Exome sequencing identifies the cause of a mendelian disorder. *Nat. Genet.* **42**, 30–35 (2010).

4. Ng, S. B. *et al.* Exome sequencing identifies *MLL2* mutations as a cause of Kabuki syndrome. *Nat. Genet.* **42**, 790–793 (2010).
5. Jagadeesh, K. A. *et al.* M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nat. Genet.* **48**, 1581–1586 (2016).
6. Ng, P. C. & Henikoff, S. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids Res.* **31**, 3812–3814 (2003).
7. Adzhubei, I. A. *et al.* A method and server for predicting damaging missense mutations. *Nat. Methods* **7**, 248–249 (2010).
8. Rehm, H. L. *et al.* ACMG clinical laboratory standards for next-generation sequencing. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **15**, 733–747 (2013).
9. Einarisdottir, E. *et al.* Identification of *NCAN* as a candidate gene for developmental dyslexia. *Sci. Rep.* **7**, 9294 (2017).
10. Bodian, D. L. *et al.* Genomic analysis of an infant with intractable diarrhea and dilated cardiomyopathy. *Mol. Case Stud.* **3**, pii: a002055, (2017).
11. Cuajungco, M. P. *et al.* Tissue-specific reduction in splicing efficiency of *IKBKAP* due to the major mutation associated with familial dysautonomia. *Am. J. Hum. Genet.* **72**, 749–758 (2003).
12. Wong, J. J.-L. *et al.* Orchestrated intron retention regulates normal granulocyte differentiation. *Cell* **154**, 583–595 (2013).
13. Marquez, Y., Höpfner, M., Ayatollahi, Z., Barta, A. & Kalyna, M. Unmasking alternative splicing inside protein-coding exons defines exitrons and their role in proteome plasticity. *Genome Res.* **25**: 995-1007 (2015).

14. Eom, T. *et al.* NOVA-dependent regulation of cryptic NMD exons controls synaptic protein levels after seizure. *eLife* **2**, e00178 (2013).
15. Faustino, N. A. & Cooper, T. A. Pre-mRNA splicing and human disease. *Genes Dev.* **17**, 419–437 (2003).
16. Scotti, M. M. & Swanson, M. S. RNA mis-splicing in disease. *Nat. Rev. Genet.* **17**, 19–32 (2016).
17. Sibley, C.R., Blazquez, L., Ule, J. Lessons from non-canonical splicing. *Nat. Rev. Genet.* **17**, 407–421 (2016).
18. Samuels, D. C. *et al.* Finding the lost treasures in exome sequencing data. *Trends Genet.* **29**, 593–599 (2013).
19. Stenson, P. D. *et al.* The Human Gene Mutation Database: towards a comprehensive repository of inherited mutation data for medical research, genetic diagnosis and next-generation sequencing studies. *Hum. Genet.* **136**, 665–677 (2017).
20. Stenson, P. D. *et al.* The Human Gene Mutation Database (HGMD) and its exploitation in the fields of personalized genomics and molecular evolution. *Curr. Protoc. Bioinforma. Ed. Board Andreas Baxevanis Al Chapter 1*, Unit1.13 (2012).
21. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
22. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
23. Ionita-Laza, I., McCallum, K., Xu, B. & Buxbaum, J. D. A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nat. Genet.* **48**, 214–220 (2016).

24. Huang, Y.-F., Gulko, B. & Siepel, A. Fast, scalable prediction of deleterious noncoding variants from functional and population genomic data. *Nat. Genet.* **49**, 618–624 (2017).
25. Barash, Y. *et al.* Deciphering the splicing code. *Nature* **465**, 53–59 (2010).
26. Zhang, C. *et al.* Integrative modeling defines the Nova splicing-regulatory network and its combinatorial controls. *Science* **329**, 439–443 (2010).
27. Xiong, H. Y. *et al.* The human splicing code reveals new insights into the genetic determinants of disease. *Science* **347**, 1254806 (2015).
28. Barbosa-Morais, N. L. *et al.* The evolutionary landscape of alternative splicing in vertebrate species. *Science* **338**, 1587–1593 (2012).
29. Cummings, B. B. *et al.* Improving genetic diagnosis in Mendelian disease with transcriptome sequencing. *Sci. Transl. Med.* **9**, eaal5209 (2017).
30. Consortium, T. Gte. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
31. Meienberg, J. *et al.* New insights into the performance of human whole-exome capture platforms. *Nucleic Acids Res.* **43**, e76 (2015).
32. Yeo, G. & Burge, C. B. Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J. Comput. Biol. J. Comput. Mol. Cell Biol.* **11**, 377–394 (2004).
33. Paggi, J. M. & Bejerano, G. A sequence-based, deep learning model accurately predicts RNA splicing branchpoints. *bioRxiv* 185868 (2017). doi:10.1101/185868
34. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).

35. Takata, A., Matsumoto, N. & Kato, T. Genome-wide identification of splicing QTLs in the human brain and their enrichment among schizophrenia-associated loci. *Nat. Commun.* **8**, ncomms14519 (2017).
36. Zhang, X. *et al.* Identification of common genetic variants controlling transcript isoform variation in human whole blood. *Nat. Genet.* **47**, 345–352 (2015).
37. Petrovski, S., Wang, Q., Heinzen, E. L., Allen, A. S. & Goldstein, D. B. Genic intolerance to functional variation and the interpretation of personal genomes. *PLOS Genet.* **9**, e1003709 (2013).
38. Huang, N., Lee, I., Marcotte, E. M. & Hurles, M. E. Characterising and predicting haploinsufficiency in the human genome. *PLOS Genet.* **6**, e1001154 (2010).
39. Samocha, K. E. *et al.* Regional missense constraint improves variant deleteriousness prediction. *bioRxiv* 148353 (2017). doi:10.1101/148353
40. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010).
41. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
42. Mort, M. *et al.* MutPred Splice: machine learning-based prediction of exonic variants that disrupt splicing. *Genome Biol.* **15**, R19 (2014).
43. Friedman, J. H. Greedy function approximation: A gradient boosting machine. *Ann. Stat.* **29**, 1189–1232 (2001).
44. Kremer, L. S. *et al.* Genetic diagnosis of Mendelian disorders via RNA sequencing. *Nat. Commun.* **8**, 15824 (2017).

45. Wang, E. T. *et al.* Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**, 470–476 (2008).
46. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An expanded view of complex traits: From Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).
47. Heid, C. A., Stevens, J., Livak, K. J. & Williams, P. M. Real time quantitative PCR. *Genome Res.* **6**, 986–994 (1996).
48. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164–e164 (2010).
49. Cunningham, F. *et al.* Ensembl 2015. *Nucleic Acids Res.* **43**, D662–669 (2015).

Figures and Tables

Table 1

	Method		
Region	CADD	LINSIGHT,EIGEN,SPIDEX	MutPred Splice
	Author's Thresholds		
	≥20	N/A	≥0.6
3' intronic	99%	N/A	N/A
3' core	4%	N/A	N/A
exonic	98%	N/A	48%
5'core	2%	N/A	N/A
5' extended	96%	N/A	N/A
5' intronic	97%	N/A	N/A
Total	31%	N/A	48%

Table 1. Misclassification rate of existing metrics at author-provided thresholds. Averaging across all regions, the CADD and MutPred Splice author-recommended thresholds result in the misclassification of 31% and 48% of pathogenic variants, respectively. It should be noted that MutPred Splice can only be applied to synonymous variants in exonic regions. Other existing tools, like LINSIGHT, EIGEN and SPIDEX, do not provide a default threshold, rendering them difficult to use for classification purposes.

Table 2

region	CADD	LINSIGHT	EIGEN	SPIDEX	MutPred Splice	S-CAP
3' intronic	≥ 0.709	≥ 0.048	≥ 3.57	≤ -1.45	N/A	≥ 0.005
3' core	≥ 21.50	≥ 0.767	≥ 7.25	≤ 1.14	N/A	Dom.: ≥ 0.031
						Rec.: ≥ 0.144
Exonic	≥ 0.061	≥ 0.086	≥ 3.78	≤ -2.31	≥ 0.090	≥ 0.012
5' core	≥ 22.6	≥ 0.795	≥ 8.381	≤ 1.65	N/A	Dom.: ≥ 0.032
						Rec.: ≥ 0.357
5' extended	≥ 7.423	≥ 0.211	≥ 13.80	≤ -0.910	N/A	≥ 0.003
5' intronic	≥ 0.852	≥ 0.047	≥ 0.847	≤ -1.636	N/A	≥ 0.004

Table 2. High sensitivity thresholds after recalibrating each method. We retuned each existing method for each of the six regions by finding the smallest threshold that resulted in the correct classification of 95% of pathogenic variants from that region. S-CAP thresholds for all regions (including the dominant and recessive modes) are included in the last column. By definition, with the high sensitivity thresholds each method will misclassify at most 5% of known pathogenic mutations, but their effectiveness at correctly classifying benign variants and reducing the number of patient VUS will vary greatly.

Table 3

Patient_ID	Disease	chr:pos (hg19)	Zygosity	Region	SCAP (%-ile)	SPIDEX (%-ile)	CADD (%-ile)	LINSIGHT (%-ile)	EIGEN (%-ile)
DDDP102313	Koolen-de-vries syndrome	17:44144914 C>T	Heterozygous	5' Core	0.081 (63.3)	-14.532(52.6)	25.8 (77.5)	0.978 (65.6)	15.932 (49.0)
DDDP100243	Claes-Jensen type mental retardation	X:53245380 C>A	Hemizygous	3' Core	0.919 (99.1)	-1.305 (22.1)	24.8 (65.02)	0.822 (23.2)	N/A
DDDP110794	Mental Retardation	3:71037144 C>T	Heterozygous	5' Core	0.217 (88.2)	-23.241(69.6)	26.4 (86.9)	0.988 (90.8)	20.234 (73.9)
DDDP111322	Nephrotic Syndrome	19:36333453 C>T	Compound Heterozygous	3' Core	0.629 (38.5)	-1.67 (24.2)	25 (69.8)	0.97 (54.8)	14.702 (40.2)
DDDP102111	Epileptic encephalopathy	2:166165305 G>A	Heterozygous	5' Core	0.348 (93.6)	-32.609 (82.8)	25 (59.8)	0.988 (90.8)	26.022 (87.7)
DDDP108825	Cohen Syndrome	8:100729602 G>A	Compound Heterozygous	5' Core	0.879 (94.2)	-7.189 (35.5)	25.1 (62.2)	0.982 (79.2)	25.405 (85.5)
DDDP100128	Sotos Syndrome	5:176673677 A>G	Heterozygous	3' Core	0.137 (81.4)	-12.574 (61.1)	22.9 (24.8)	0.961 (94.5)	9.876 (18.6)
DDDP111152	Ehlers-Danlos Syndrome	2:189873948 G>A	Heterozygous	5' Core	0.777 (98.6)	-5.594 (31.1)	27.9 (97.5)	0.989 (92.9)	29.666 (90.5)
DDDP102594	Mental Retardation	X:135095506 A>G	Hemizygous	3' Core	0.758 (98.6)	-2.955 (31.0)	24.3 (52.9)	0.874 (29.6)	N/A
DDDP111486	Noonan Syndrome	11:11914887 4 A>T	Heterozygous	3' Core	0.539 (97.6)	-28.352 (89.7)	24.9 (67.4)	0.982 (79.3)	22.809 (80.2)
DDDP100281	Epileptic encephalopathy	16:56370773 G>A	Heterozygous	5' Core	0.151 (82.1)	-7.739 (37.0)	26.6 (89.2)	0.987 (88.5)	24.055 (84.9)
C11	Congenital fiber-type disproportion	19:38958362 C>T	Heterozygous	Exonic	0.657 (97.8)	-0.366 (89.6)	19.430 (97.6)	N/A	7.879 (19.8)
C1	Dystroglycanopathy	1:46655129 C>A	Compound Heterozygous	5' Core	0.817 (77.03)	-52.028 (95.9)	21.500 (14.3)	0.992 (96.8)	4.976 (7.63)
		1:46660532 G>A		Exonic	0.095 (83.0)	-1.556 (96.4)	12.050 (79.2)	N/A	21.658 (76.5)
E2	Nemaline myopathy	2:152520057 C>T	Compound Heterozygous	5' Extended	0.147 (98.9)	-11.190 (85.4)	19.680 (97.9)	0.981 (79.2)	20.272 (65.0)
		2:152544805 C>T		5' Core	0.915 (99.1)	-2.970 (22.6)	27.200 (94.1)	0.990 (94.7)	20.384 (74.5)

Table 3. Causative variants in each patient and pathogenicity predictions. Each row describes a single patient, their underlying disease, causative variant and its zygosity, the region of the associated gene in which it is located, and the score and percentile each method assigns to the variant. The percentile was computed by measuring the fraction of variants in the same region with a score less than the score assigned to the causative variant. The 99th percentile denotes that 99% of variants in that region have a score less pathogenic than the one we are observing thereby indicating that the variant is considered to be highly pathogenic. Dark and light green filled entries have the highest and second highest percentile scores for the given variant,

respectively. Red entries highlight patients where the causative variant would have been classified as benign using the author-recommended thresholds.

Figure 1

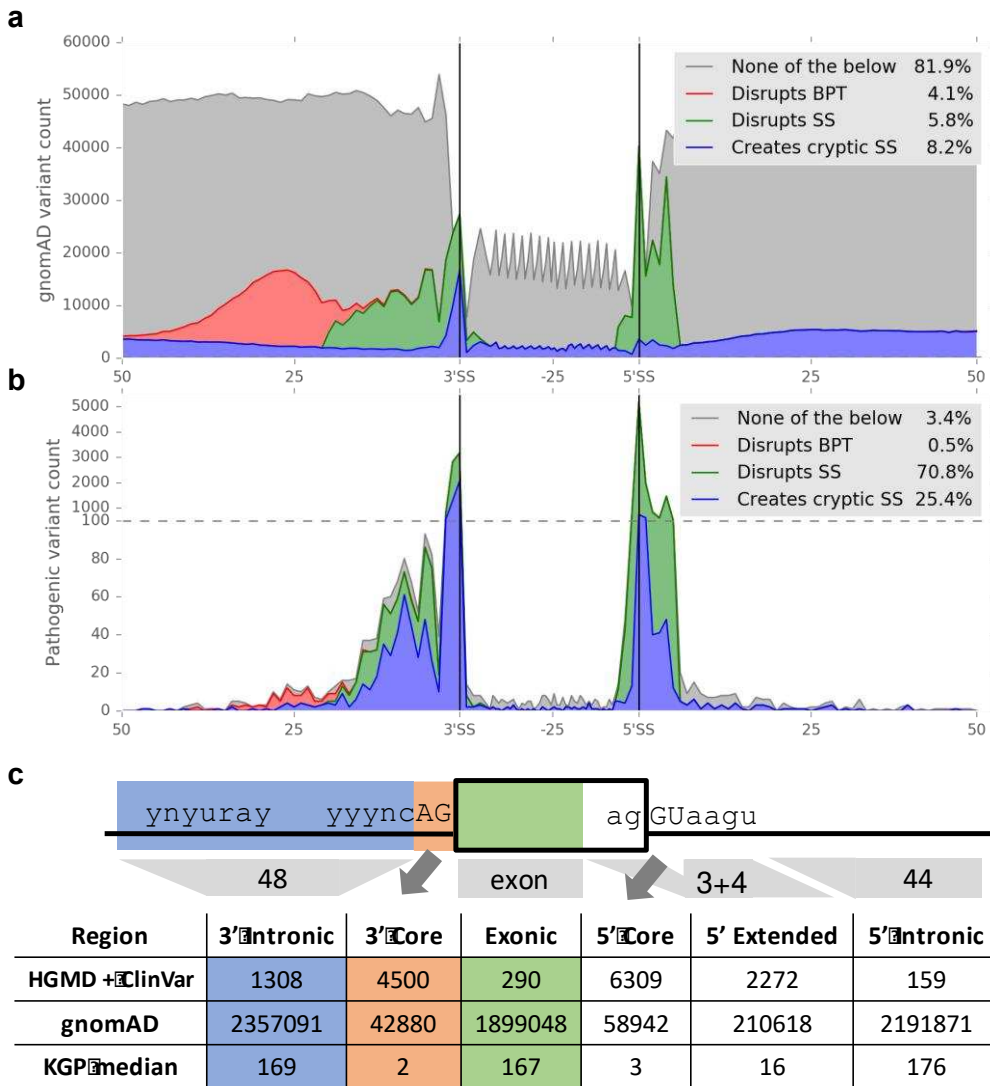


Figure 1. Distribution of rare, noncoding variants in the splicing region. We built two sets of rare variants near exons that have no effect on the annotated coding sequence: (A) a set of likely benign variants from the Genome Aggregation Database (gnomAD) and (B) a set of putatively pathogenic variants from the Human Gene Mutation Database (HGMD) and ClinVar. We developed a simple model to ascertain the effect of these variants on splicing, using MaxEntScan to assess if a variant created a cryptic splice site or disrupted the reference splice site, and LaBranchoR *in silico* mutagenesis scores to assess if a variant disrupted a branchpoint. In (A) and (B), we plot the aggregate counts of each variant set as a function of position relative to the nearest splice site, colored by their putative effect. Nearly 97% of pathogenic variants in the splicing region [GILL: what do you mean by the 'splicing region'?] are predicted to have an effect [GILL: Do you mean a deleterious effect or any effect?] on splicing, as compared to only 18% of benign variants. (C) We split the variants into different regions with largely homogenous function. The majority of known pathogenic variants are found in the 3' and 5' core splice site regions, whereas the majority of benign variants are found in the 3' intronic, exonic and 5' intronic

regions. In a typical individual, the distribution of variants is similar to the distribution of gnomAD benign variants.

Figure 2

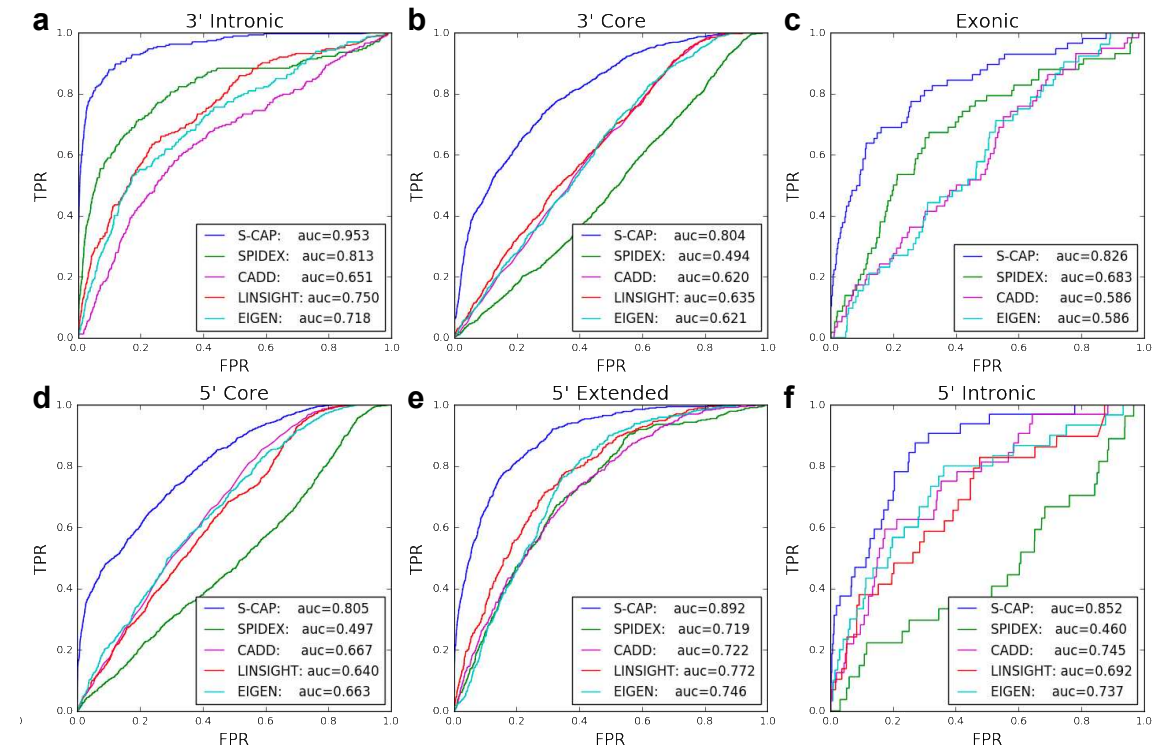


Figure 2. Overall performance per region for splicing pathogenicity classification. For each method, 1000 threshold points were determined by evenly spanning the range from the minimum to the maximum score observed for the method. A true positive rate and false positive rate were determined for each threshold value and used to build the receiver operating characteristic (ROC) curve. S-CAP achieves an AUC of 0.953 in the 3' intronic region (A), an AUC of 0.804 in the 3' core sites (B), an AUC of 0.826 in the exonic region (C), an AUC of 0.805 in the 5' core sites (D), an AUC of 0.892 in the 5' extended region (E) and an AUC of 0.852 in the 5' intronic region (F). S-CAP outperforms existing metrics in all regions and no existing method consistently outperforms the rest [GILL:.....whereas none of the existing method consistently outperformed the others?].

Figure 3.

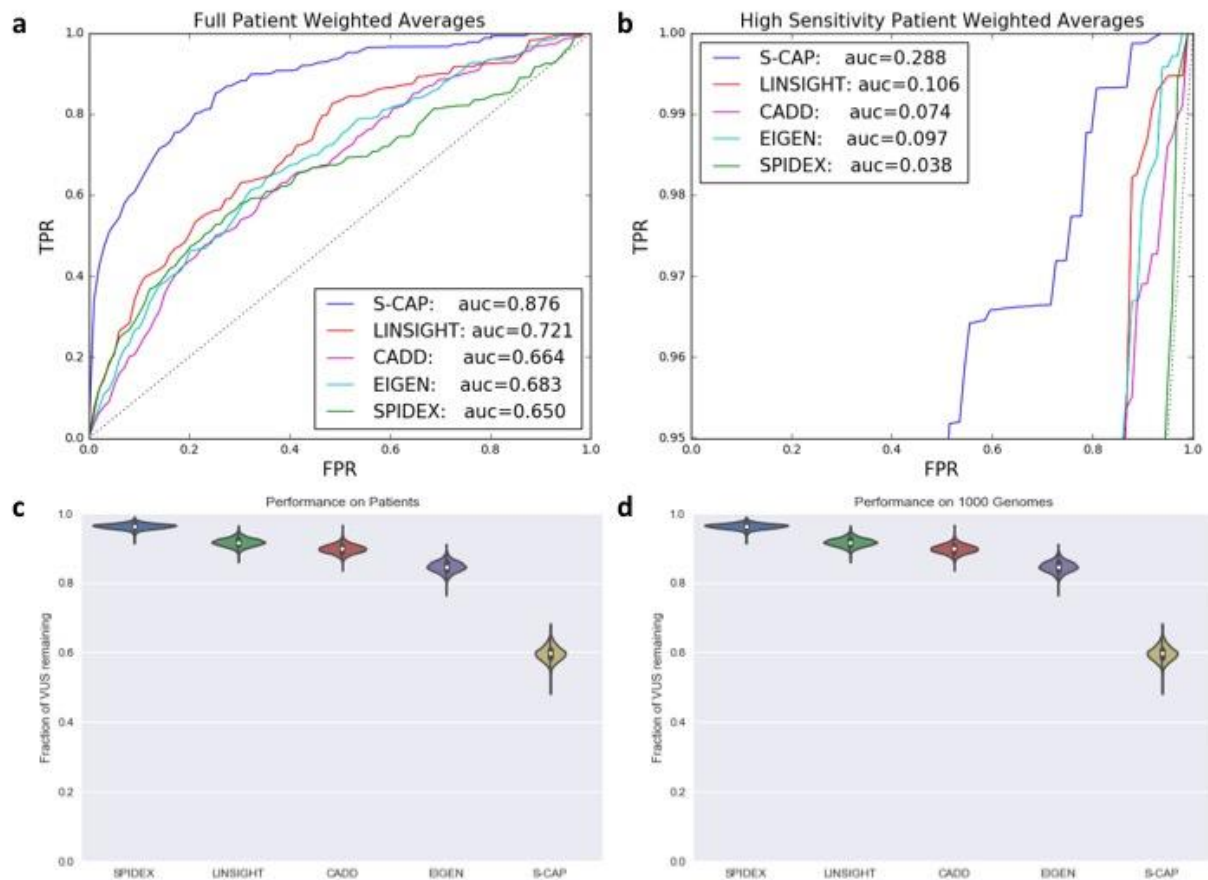


Figure 3. Overall performance on patient data. Since we have defined separate models for 6 different regions and separate the 2 core regions into dominant and recessive, there are 8 different AUC curves for each method. We take a weighted sum of each AUC based on the distribution of variants seen in a typical individual to form an overall receiver operating characteristic (ROC) curve representative of the overall performance expected on patients. (A) S-CAP achieves an AUC of 0.876 and (B) an hsr-AUC of 0.288. (C) S-CAP reduces the number of splicing related variants of uncertain significance (VUS) from patient exomes by 40% while maintaining the pathogenic variants with 95% sensitivity. At the same sensitivity requirement, existing methods reduce the VUS by at most 15%. (D) We observe a similar reduction in VUS overall (n=2054) Thousand Genomes Project individuals, which conceptually only differ from Mendelian disease patients by up to 2 mutations.

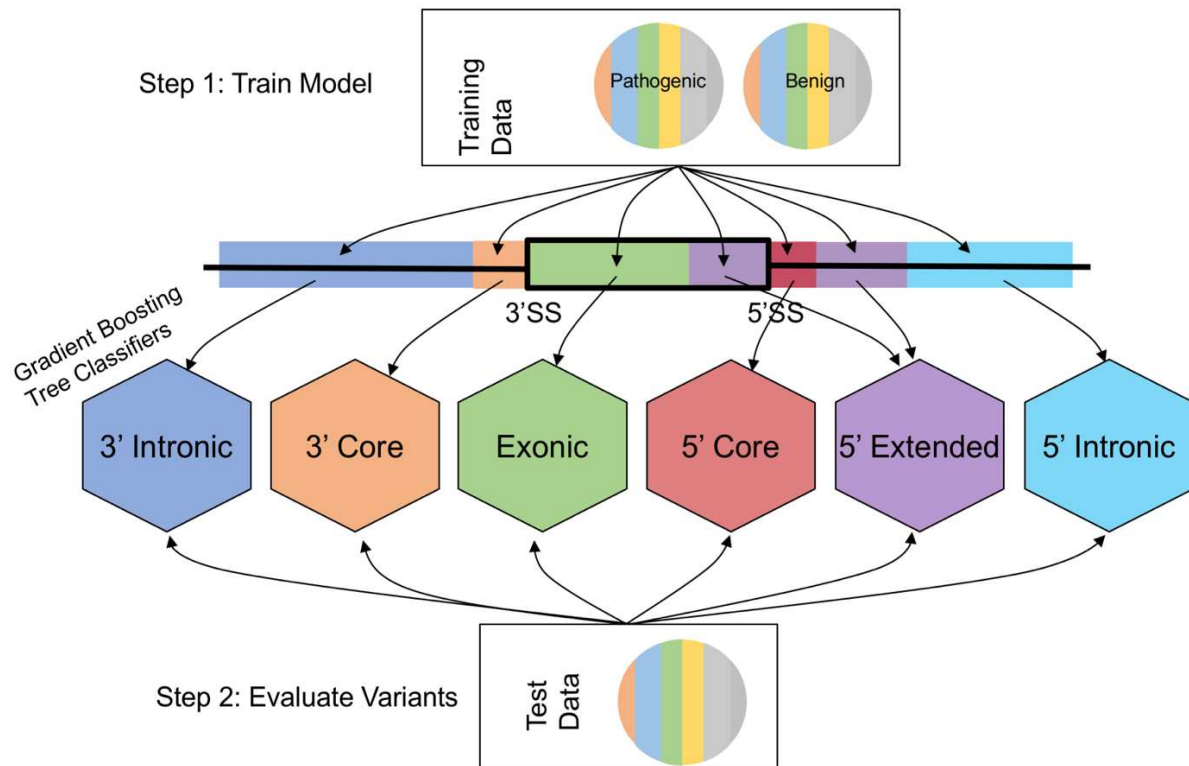
Supplementary Figures

Supplementary Table 1

Feature Category	Features	Description
Chromosome Level	X, Y or not XY	3 binary features to indicate whether a variant is found on the autosomes or the X or Y chromosome
Gene Level	RVIS	Measure of gene mutability. Observed v. Expected variant abundance
	pLI	Measure of gene loss of function tolerance.
	haploinsufficiency score	Recessive or dominant gene inheritance
Exon Level	# rare 3' core variants	rare variants observed in 3' core region
	# rare 5' core variants	rare variants observed in 5' core region
	# common 3' core variants	common variants observed in 3' core region
	# common 5' core variants	common variants observed in 5' core region
	# rare 3' extended variants	rare variants observed in 3' extended region
	# rare 5' extended variants	rare variants observed in 5' extended region
	# common 3' extended variants	common variants observed in 3' extended region
	# common 5' extended variants	common variants observed in 5' extended region
	exon % identity	Measure the % identity of hg19 exon with each of the 99 species in the multiz100way. Take top 6 PCA components after fitting a PCA to all coding exons.
	exon length	Number of bases in the exon
	exon length % 3	Number of bases modulo 3
	MPC	Regional mutational constraint score
Variant Level	Spectrum Kernel	Count of all 3-mers introduced and removed by mutation
	MaxEntScan	The difference in motif match
	SPIDEX	Predicted $\Delta\psi$ (change in exon expression)
	distance to 5' SS	Number of bases to the 5' splice site
	distance to 3' SS	Number of bases to the 3' splice site
	CADD	Functional data SVM-based classifier
	LINSIGHT	Conservation based model to identify variants under negative selection.
	PhyloP	Base-pair conservation across primates, mammals and vertebrates
	PhastCons	Regional conservation across primates, mammals and vertebrates
3' Intronic	LaBranchoR	Sequence based deep learning branchpoint prediction
3' Intronic, 3' Core, Exonic	3' cryptic splice site creation terms	MaxEntScan based feature to measure cryptic splice creation near the 3' side
Exonic, 5' Core, 5' Extended, 5' Intronic	5' cryptic splice site creation terms	MaxEntScan based feature to measure cryptic splice creation near the 5' side
3' Core, 5' Core	Zygosity	Indicates whether the variant is seen in a heterozygous or homozygous state.

Supplementary Table 1. Description of features used to build S-CAP. The chromosome, gene, exon and variant level features were used in all models for all regions. There was an additional set of features that was specific to certain regions. These are enumerated in the table below the variant level features section.

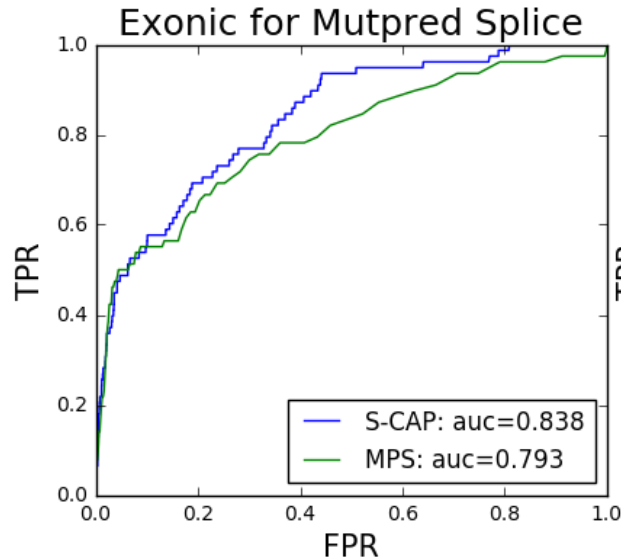
Supplementary Figure 1



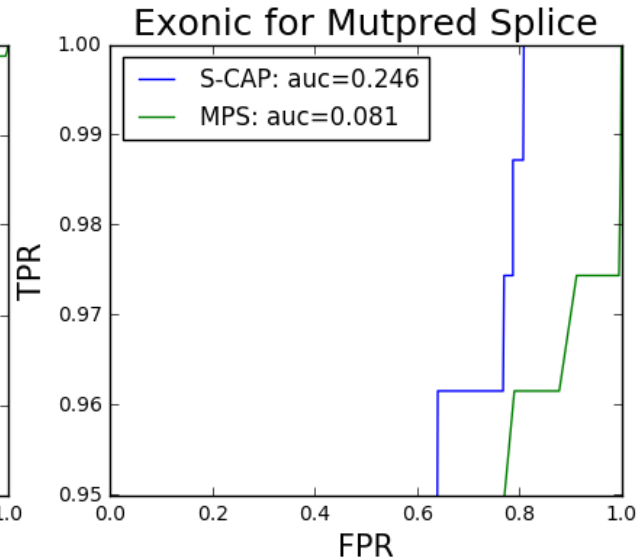
Supplementary Figure 1. Framework for training and evaluating 6 pathogenicity models. The splicing region is split into 6 independent regions as defined in **Fig. 1c**, and a separate model is trained for variants residing in each region. For the 5' and 3' core regions, we additionally defined separate dominant and recessive classifiers (not shown) leading to a total of 8 models. Given a set of variants to be scored, we calculate the S-CAP score for each variant by using the corresponding model associated with the region where the variant is found.

Supplementary Figure 2.

A

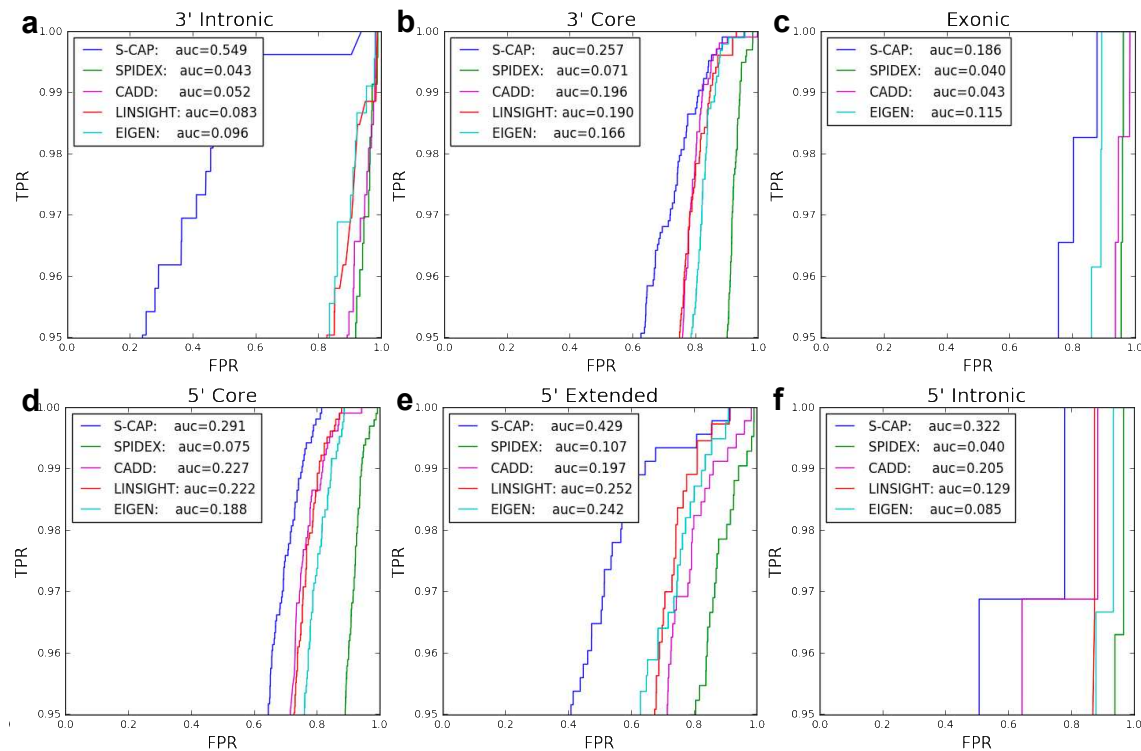


B



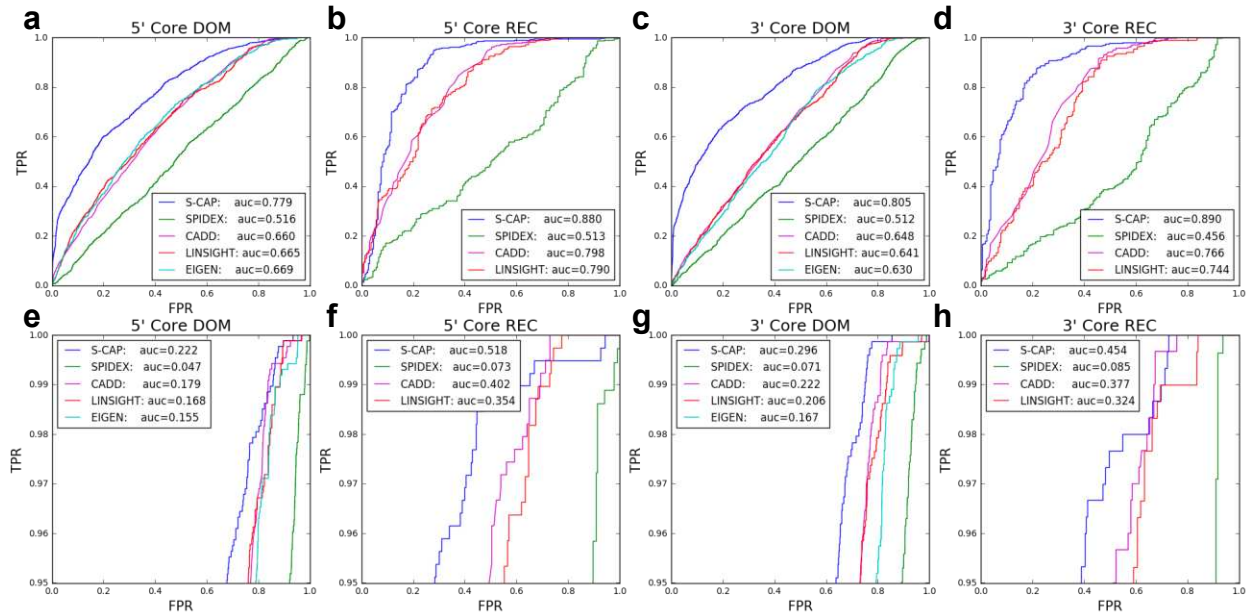
Supplementary Figure 2. Performance of S-CAP compared to MutPred Splice. MutPred Splice is a computational method for predicting the pathogenicity of exonic synonymous variants. MutPred Splice was trained by its authors using a subset of the pathogenic data used to train/test S-CAP. As a result, we need to independently test MutPred Splice on a set of variants that was not used in its training. This test set comprises rare synonymous variants from HGMD added to the database in 2013 or later. On this set, S-CAP achieves an AUC of 0.838 while MutPred Splice achieves an AUC of 0.793. S-CAP performs especially well in the high sensitivity domain with a hsr-AUC of 0.246 as compared to 0.081 for MutPred Splice.

Supplementary Figure 3.



Supplementary Figure 3. S-CAP performance in the high sensitivity region. The hsr-AUC curve is formed by subsetting the overall AUC to just the region where pathogenic variants are correctly classified over 95% of the time. An hsr-AUC curve is calculated for each of the regions as defined in **Fig. 1C**. S-CAP achieves an hsr-AUC of 0.549 in the 3' intronic region (A), an hsr-AUC of 0.257 in the 3' core sites (B), an hsr-AUC of 0.186 in the exonic region (C), 0.291 in the 5' core sites (D), 0.429 in the 5' extended region (E) and 0.322 in the 5' intronic region (F). S-CAP outperforms existing metrics in all regions in the high sensitivity domain whereas none of the existing methods consistently outperforms the others [GILL: Is this what you want to say?].

Supplementary Figure 4



Supplementary Figure 4. Performance on recessive and dominant classes. The distribution of the underlying features is dramatically different for dominant and recessive variants. This results in a big difference in performance when classifying recessive and benign variants in the core splice site regions. S-CAP achieves an AUC (A) of 0.779 on dominant tagged variants and (B) of 0.880 on recessive tagged variants in the 5' core region. There is a similar performance difference in the 3' core region where S-CAP achieves an AUC (C) of 0.805 on dominant tagged variants (D) and of 0.890 on recessive tagged variants. In the high sensitivity region, S-CAP achieves an hsr-AUC (E, F) of 0.222 on dominant tagged variants and of 0.518 on recessive tagged variants in the 5' core region (G, H) and of 0.296 on dominant tagged variants and of 0.454 of recessive tagged variants in the 3' core region.